

A subgroup-aware scoring approach to the study of effect modification in observational studies

Yijun Fan^a, Dylan S. Small^b

^a*Graduate Group of Applied Mathematics and Computational Science, University of Pennsylvania, Philadelphia, 19104, PA, USA*

^b*Department of Statistics and Data Science, University of Pennsylvania, Philadelphia, 19104, PA, USA*

Abstract

Effect modification means the size of a treatment effect varies with an observed covariate. Generally speaking, a larger treatment effect with more stable error terms is less sensitive to bias. Thus, we might be able to conclude that a study is less sensitive to unmeasured bias by using these subgroups experiencing larger treatment effects. Lee et al. (2018) proposed the submax method that leverages the joint distribution of test statistics from subgroups to draw a firmer conclusion if effect modification occurs. However, one version of the submax method uses M-statistics as the test statistics and is implemented in the R package `submax` (Rosenbaum, 2017). The scaling factor in the M-statistics is computed using all observations combined across subgroups. We show that this combining can confuse effect modification with outliers. We propose a novel group M-statistic that scores the matched pairs in each subgroup to tackle the issue. We examine our novel scoring strategy in extensive settings to show the superior performance. The proposed method is applied to an observational study of the effect of a malaria prevention treatment in West Africa.

Keywords: Sensitivity analysis, Power of a sensitivity analysis, Effect modification, Treatment heterogeneity, Subgroup analysis

1. Introduction

1.1. Effect modification and sensitivity analysis

In an observational study, subjects are not randomly assigned to treatment or control, so the treated and control groups might differ in terms of both measured and unmea-

sured covariates. To make subjects in the treated group and control group comparable, adjustments (e.g. matching) are used to address the concern of measured covariates, but there is still typically concern about unmeasured covariates. Then, a sensitivity analysis asks how large the magnitude of bias from some unmeasured covariates would need to be to explain away the qualitative conclusion of a study based on the assumption of no unmeasured confounding (Rosenbaum, 2002).

Effect modification occurs when the size of a treatment effect varies depending on a measured covariate. In general, larger treatment effects are less sensitive to unmeasured confounding, which suggests that effect modification can play a role in reducing sensitivity to unmeasured covariates (Hsu et al., 2013). If some subgroups experience larger effects and we make use of those subgroups appropriately, then we may be able to report less sensitivity to unmeasured confounding. The submax method, proposed by Lee et al. (2018), splits the population into certain overlapped subgroups and uses the joint distribution of these correlated test statistics from the subgroups to study the effect modification in observational studies. For example, a researcher may want to make a robust inference for subgroup analyses by using the submax method based on Huber M-statistics following the suggestion of Maritz (1979). However, although the submax method considers the subgroups, the scores entering the chosen test statistic (see section 2.1 and 3.1) used by the submax method do not account for such subgroup structures. In this article, we would show that such practice could substantially deteriorate the sensitivity analysis in some cases because effect modification is confused with outliers. We illustrate this point by using M-statistics (Maritz, 1979) and propose a subgroup-aware scoring approach.

The motivating example of malaria control in West Africa is introduced in section 1.2. The notation and the submax method is reviewed in section 2. We introduce the novel scoring approach in section 3 and apply the proposed method to the malaria control example in section 4. A simulation is provided in section 5 to evaluate this approach in extensive settings.

1.2. Motivating example: control of malaria in West Africa

The World Health Organization, with the Nigerian government, worked on comparing the effectiveness of different treatment methods of controlling malaria in West Africa

(Molineaux and Gramiccia, 1980). The treatment of interest in our study is spraying with propoxur, an insecticide, and mass administration of a drug called sulfalene-pyrimethamine. The outcome measurement documents the frequency of *Plasmodium falciparum*, a protozoan parasite that causes malaria, in blood samples. A series of blood samples were collected across a period of time and we used the four surveys immediately before the treatment and four surveys immediately after the treatment to compute a score. Specifically, any individual included was required to have at least two surveys before and after treatment; then these surveys were combined by using Huber’s M-estimate as a trimmed mean for pretreatment and posttreatment summaries. We are interested in whether such treatment would cause lower *Plasmodium falciparum* frequency. Prior to using the outcome information other than quality control, we did pair matching for age and gender, as in Hsu et al. (2013), which resulted in $I = 1560$ matched pairs.

2. Setup of sensitivity analysis for subgroup comparisons

2.1. Notation

Suppose we are interested in L binary covariates as potential effect modifiers and thus $G = 2^L$ non-overlapping interaction subgroups are defined. Specifically, in our example for malaria control, we are interested in $L = 2$ covariates (i.e. age and gender), and consequently $G = 2^2 = 4$ independent interaction groups are defined. Consider here we have I_g matched sets in each group g , $1 \leq g \leq G$, and in each matched set gi , $1 \leq i \leq I_g$, there is one treated subject j and $n_{gi} - 1$ non-treated subjects. For simplicity, we only consider the case of matched pairs here, that is $n_{gi} = 2$ and we have the measured covariates $x_{gi1} = x_{gi2}$. In addition to the measured covariates, there may be an unmeasured covariate $u_{gi1} \neq u_{gi2}$. Denote Z_{gij} as an indicator of treatment assignment, so we have $1 = Z_{gi1} + Z_{gi2}$. We denote r_{Tgij} and r_{Cgij} as the responses of j^{th} individual in the i^{th} matched pair of group g . Following the potential outcome framework, we have the exhibited response for individual gij , that is $R_{gij} = Z_{gij}r_{Tgij} + (1 - Z_{gij})r_{Cgij}$. Write $\mathcal{F} = \{(r_{Tgij}, r_{Cgij}, x_{gij}, u_{gij}), g = 1, \dots, G, i = 1, \dots, I_g, j = 1, 2\}$. Denote \mathcal{Z} as the set containing $|\mathcal{Z}| = \prod_{g=1}^G 2^{I_g}$ possible values \mathbf{z} of all the possible treatment assignments $\mathbf{Z} = (Z_{111}, Z_{112}, \dots, Z_{G,I_G,1}, Z_{G,I_G,2})^T$, so $\mathbf{z} \in \mathcal{Z}$ if $z_{gij} \in \{0, 1\}$ and $1 = z_{gi1} + z_{gi2}$ for

each g and i . Conditioning on the event $\mathbf{Z} \in \mathcal{Z}$ is abbreviated as conditioning on \mathcal{Z} , and we have $\Pr(\mathbf{Z} = \mathbf{z} \mid \mathcal{F}, \mathcal{Z}) = |\mathcal{Z}|^{-1}$ for each $\mathbf{z} \in \mathcal{Z}$.

We are interested in testing Fisher's sharp null hypothesis of no treatment effect $H_0 : r_{Tgij} = r_{Cgij}, g = 1, \dots, G, i = 1, \dots, I_g, j = 1, 2$. Under the null, we have the observed treated-minus-control difference $D_{gi} = (Z_{gi1} - Z_{gi2})(R_{gi1} - R_{gi2}) = (Z_{gi1} - Z_{gi2})(r_{Cgi1} - r_{Cgi2}) = \epsilon_{gi}$ of matched pair gi . If treatments were assigned randomly, the randomization creates the exact null distribution of the statistic $T = \sum_{g=1}^G T_g = \sum_{g=1}^G \sum_{i=1}^{I_g} \sum_{j=1}^2 Z_{gij} q_{gij}$ with some properly chosen score function q_{gij} , since now the only randomness is from the treatment assignment \mathbf{Z} as we condition on \mathcal{F} and \mathcal{Z} . The score q_{gij} may depend on R_{gij} or x_{gij} . For example, it replicates Maritz's version of Huber M-statistic (Maritz, 1979) if we use the trimming version of the score (see section 3.1) and permutational t-test if we use the raw data.

However, in an observational study, treatment and control are not assigned at random. Formally, in Rosenbaum's sensitivity model (Rosenbaum, 1987, 2002), a sensitivity parameter $\Gamma \geq 1$ is introduced. Denote $\pi_{gij} = \Pr(Z_{gij} = 1 \mid \mathcal{F})$, and we assume

$$\frac{1}{\Gamma} \leq \frac{\pi_{gij}(1 - \pi_{g'i'j'})}{\pi_{g'i'j'}(1 - \pi_{gij})} \leq \Gamma, \text{ whenever } \mathbf{x}_{gij} = \mathbf{x}_{g'i'j'};$$

Γ is the maximum odds ratio of getting treatment for two units with the same measured covariates but a different unmeasured covariate. $\Gamma = 1$ is a randomized experiment; larger Γ 's allow for more unmeasured confounding. Given a Γ , we obtain a range of possible significance levels (i.e. P -values) under this constraint on the treatment assignment. Based on the asymptotic separability technique (Gastwirth et al., 2000), we could obtain an approximation to the upper bound for P -values. Formally, denote $\mu_{\Gamma g}$ as the maximum expectation of T_g and $\nu_{\Gamma g}$ as the maximum variance of T_g that achieves the maximum expectation $\mu_{\Gamma g}$ under this constraint. When $\min(I_g) \rightarrow \infty$ together with some regularity conditions, we have the following approximated P -value upper bound

$$1 - \Phi \left\{ \left(\sum_{g=1}^G T_g - \mu_{\Gamma g} \right) / \sqrt{\sum_{g=1}^G \nu_{\Gamma g}} \right\},$$

where Φ is the cumulative distribution function of the standard normal distribution. Moreover, now that we have $\min(I_g) \rightarrow \infty$ and G groups are independent, it holds that

the joint distribution of G statistics $(T_g - \mu_{\Gamma g})/\nu_{\Gamma g}^{1/2}$ is a G -dimensional standard normal distribution.

2.2. The submax method: joint bounds of subgroup comparisons

Now that we wish to jointly evaluate multiple subgroups, the submax method (Lee et al., 2018) uses the maximum standard deviates of multiple subgroups as the test statistic, whose distribution under the null could be derived explicitly in the large sample as $\min(I_g) \rightarrow \infty$. However, instead of conducting $G = 2^L$ independent tests for these finest groups, the submax method only conducts $2L$ subgroup tests plus one overall test, resulting in $K = 2L + 1$ correlated tests in total. In the running example, the submax method would conduct $K = 2 \times 2 + 1 = 5$ correlated tests for all the subjects, subjects with age less than 10, subjects with age larger than 10, female, and male respectively. Formally, suppose we are interested in making $K = 2L + 1$ comparisons, and we define the $K \times G$ matrix \mathbf{C} such that each row $\mathbf{c}_{k*} = (c_{k1}, \dots, c_{kG})^T$, $1 \leq k \leq K$ encodes whether a group g , $1 \leq g \leq G$, is included in the k^{th} comparison. Then, for this comparison, we have the score $S_k = \sum_{g=1}^G c_{kg} T_g$. Write $\boldsymbol{\mu}_{\Gamma} = (\mu_{\Gamma 1}, \dots, \mu_{\Gamma G})^T$ and $\mathbf{V}_{\Gamma} = \text{diag}\{\nu_{\Gamma 1} \dots \nu_{\Gamma G}\}$. Thus, we have the joint distribution of (S_1, \dots, S_K) as $N(\mathbf{C}\boldsymbol{\mu}_{\Gamma}, \mathbf{C}\mathbf{V}_{\Gamma}\mathbf{C}^T)$. Denote $\boldsymbol{\theta}_{\Gamma} = \mathbf{C}\boldsymbol{\mu}_{\Gamma}$ and $\boldsymbol{\Sigma}_{\Gamma} = \mathbf{C}\mathbf{V}_{\Gamma}\mathbf{C}^T$ while we write $\theta_{\Gamma k}$ as the k^{th} coordinate of $\boldsymbol{\theta}_{\Gamma}$ and $\sigma_{\Gamma k}$ as the k^{th} diagonal element of $\boldsymbol{\Sigma}_{\Gamma}$. Finally, denote $D_{\Gamma k} = (S_k - \theta_{\Gamma k})/\sigma_{\Gamma k}$, and we have the joint distribution of $\mathbf{D}_{\Gamma} = (D_{\Gamma 1}, \dots, D_{\Gamma K})^T$ as $N(0, \boldsymbol{\rho}_{\Gamma})$, where $\boldsymbol{\rho}_{\Gamma}$ is the $K \times K$ correlation matrix whose element in i^{th} row and j^{th} column is derived from dividing $(\boldsymbol{\Sigma}_{\Gamma})_{ij}$ by $\sigma_{\Gamma i}$ and $\sigma_{\Gamma j}$. We now have the null distribution of the test statistic $D_{\Gamma \max} = \max_{1 \leq k \leq K} D_{\Gamma k}$ used in the submax method. The critical value $\kappa_{\Gamma, \alpha}$ at level α for $D_{\Gamma \max}$ solves $1 - \alpha = \Pr(D_{\Gamma \max} < \kappa_{\Gamma, \alpha})$. In particular, in the case of matched pairs, it holds that $\rho_{\Gamma} = \rho$ would not depend on Γ (it only depends on the observed data). Thus, given the input scores, the corresponding critical value $\kappa_{\Gamma, \alpha} = \kappa_{\alpha}$ does not vary with Γ in the case of matched pairs.

3. Subgroup-aware scoring approach

3.1. M-score by trimming together

Before introducing our proposed method, it might be helpful to review the conventional practice when using Maritz-Huber’s M-statistic in the submax method. Suppose we are interested in testing the Fisher’s sharp null hypothesis of no treatment effect H_0 using M-statistic, so we plug in the M-score. Denote $h_0 = \text{median}\{|D_1|, |D_2|, \dots, |D_I|\}$, and we have

$$T = \sum_{g=1}^G T_g = \sum_{g=1}^G \sum_{i=1}^{I_g} \text{sign}(D_{gi}) \psi\left(\frac{|D_{gi}|}{h_0}\right),$$

where $\psi(\cdot)$ is an odd function and $\psi(d)$ is non-negative when $d > 0$; $\text{sign}(\cdot)$ indicates the sign. Common in the practice is the trimming version of $\psi(\cdot)$ suggested by Maritz (1979) and further discussion about its application in sensitivity analysis is provided in Rosenbaum (2007). Specifically, $\psi(d) = 0$ when $|d| < a$ while $\psi(d) = d$ when $a < |d| \leq t$ and $\psi(d) = t$ when $|d| > t$. a is the so-called inner parameter (Rosenbaum, 2013) and t is the so-called trimming parameter. In this article, we follow the default setting of function `Mscorev` in the R package `submax` (Rosenbaum, 2017) where $a = 0$ and $t = 3$ when carrying out the experiments.

3.2. The proposed method

We are ready to introduce our novel scoring approach that accounts for the subgroup structures for the aforementioned trimming version of M-statistics. Instead of applying trimming for all the matched pairs simultaneously (i.e. scalar h_0 is applied for all D_{gi}), we employ this trimming procedure within each group g . Formally, denote $h_{g0} = \text{median}\{|D_{g1}|, |D_{g2}|, \dots, |D_{gI_g}|\}$, and we have the following group M-statistic that uses the subgroup structure:

$$T^{sub} = \sum_{g=1}^G T_g^{sub} = \sum_{g=1}^G \sum_{i=1}^{I_g} \text{sign}(D_{gi}) \psi\left(\frac{|D_{gi}|}{h_{g0}}\right) h_{g0}.$$

In particular, without the subgroup argument (i.e. $G = 1$), it reduces to the M-statistic applied to all the matched pairs since there will be no difference if a constant (say h_0 here) is multiplied to all the observations in a permutation test. Note that when $G > 1$, for each group g , we use the scalar h_{g0} to do the trimming and retain the relative differences

of the score sizes across different subgroups by multiplying the weight by the within-group median h_{g0} . The ψ function is the same as in section 3.1 for all the numerical results in the article.

4. Application in the malaria example

As shown in Table 1, we applied the proposed group M-statistic in section 3.2 to the malaria control example. In this example, effect modification takes place in the subgroup of young people (with age less than 10 years old), and the maximum statistic is always based on those 447 pairs of young people. It is notable that if we use the raw data (i.e. mean difference statistic) as the score, we still have evidence to reject the null hypothesis of no treatment effect up to the sensitivity value at $\Gamma = 4.1$; however, if we use the M-statistic by trimming all the matched pairs together, then we have evidence up to the sensitivity value only at $\Gamma = 2.7$. In other words, if the researcher wants to use M-statistics to make a robust argument, the reported insensitivity to bias is substantially less than using the raw data. Intuitively, since all the treated-minus-control differences are scored together, those treated-minus-control differences with large and positive values from some subgroup where effect modification happens are affected most by the trimming in the M-statistic, resulting in the unwanted smaller Γ . Our proposed group M-statistic in section 3.2 employs the trimming process within each subgroup and Table 1 shows that we are still able to reject the null hypothesis of no treatment effect with a sensitivity value at $\Gamma = 4.0$, comparable to using the raw data.

5. Simulation

In this section, we examine the power of sensitivity analysis in finite samples of the proposed method in extensive settings to show the superior performance compared with using raw data and the conventional practice of using M-statistics. We evaluate the power of sensitivity analysis as in other papers under the "favorable situation" that there is in fact no unmeasured confounding (Hsu et al., 2013; Hansen et al., 2014). However, if a researcher is in this favorable situation, it could not be told from the observed data. Then, a sensitivity analysis is performed and perhaps the best we can hope to say is

k	1	2	3	4	5	6
Subpopulation	All	Age ≤ 10	Age > 10	Female	Male	Maximum
	D_{Γ_1}	D_{Γ_2}	D_{Γ_3}	D_{Γ_4}	D_{Γ_5}	$D_{\Gamma_{\max}}$
Sample-size	1560	447	1113	766	794	
$\Gamma = 2.0$						
Mean difference	5.28	6.63	-2.30	2.38	5.02	6.63
M-statistic	2.03	4.86	-2.55	0.26	2.5	4.86
Group M-statistic	6.21	6.73	-2.89	2.74	5.82	6.73
$\Gamma = 2.7$						
Mean difference	2.52	4.74	-5.00	0.48	3.02	4.74
M-statistic	-2.10	2.31	-5.92	-2.63	-0.44	2.31
Group M-statistic	3.68	4.74	-6.63	0.97	4.01	4.74
$\Gamma = 3.0$						
Mean difference	1.57	4.11	-5.97	-0.19	2.34	4.11
M-statistic	-3.56	1.43	-7.14	-3.66	-1.48	1.43
Group M-statistic	2.81	4.07	-7.98	0.35	3.40	4.07
$\Gamma = 4.0$						
Mean difference	-1.01	2.42	-8.71	-2.01	0.50	2.42
M-statistic	-7.60	-0.96	-10.55	-6.52	-4.34	-0.96
Group M-statistic	0.48	2.29	-11.78	-1.33	1.76	2.29
$\Gamma = 4.1$						
Mean difference	-1.24	2.28	-8.95	-2.17	0.35	2.28
M-statistic	-7.95	-1.17	-10.86	-6.77	-4.58	-1.17
Group M-statistic	0.28	2.14	-12.11	-1.48	1.63	2.14

Table 1: Sensitivity analysis by using submax method with three test statistics to the malaria control data. We are interested in two binary covaries, age and gender, here, so the submax method would return five test statistics and we reject the null hypothesis if the maximum of these five statistics is larger than the corresponding critical values. The critical values for using mean difference, M-statistic, and group M-statistic are 2.20, 2.20, 2.18 respectively. The test statistics larger than the critical values are in bold.

that the study is insensitive to certain unmeasured confounding. It is in this favorable situation that we evaluate the power of sensitivity analysis for different test statistics. Now consider we have two binary covariates of interest and one of them is set to be the effect modifier. A favorable situation in the effect modification setting would be that we pre-specify the constant effect sizes τ_1 and τ_2 for different subgroups and sample the independent error terms ϵ_1 and ϵ_2 that might be different across subgroups (as we are sampling conditional on the covariate) but have a symmetric distribution centered at zero and i.i.d within a group since the treatment is assigned at random. Specifically, suppose we now have $I = 1000$ matched pairs with two binary covariates and $I_g = 250, 1 \leq g \leq 4$ in the four groups formed by the two binary covariates. Effect modification only happens with the first covariate. Specifically, the first 500 matched pairs have value 1 of the first covariate and the second 500 matched pairs have value 0, and the effect of the treatment is larger in the first 500 matched pairs.

Inspired by the distribution of treated-minus-control differences of young children and older people in the malaria control data in Figure B.1, we further consider the following sampling situations of treated-minus-control difference for $I = 1000$ independent matched pairs and repeat each sampling situation for 10,000 times to estimate the power:

(1) $D_i = 5 + 10\epsilon_{1i}, 1 \leq i \leq 500$, where ϵ_{1i} is standard normal distribution; $D_i = 0.5 + \epsilon_{2i}, 501 \leq i \leq 1000$, where ϵ_{2i} is t distribution with degree of freedom 2.

(2) $D_i = 5 + 5\epsilon_i, 1 \leq i \leq 500$; $D_i = 0.5 + 0.5\epsilon_i, 501 \leq i \leq 1000$, where ϵ_i is t distribution with degree of freedom 3.

(3) $D_i = 4 + 5\epsilon_i, 1 \leq i \leq 500$; $D_i = 0.2 + \epsilon_i, 501 \leq i \leq 1000$, where ϵ_i is standard normal distribution.

(4) $D_i = 5 + 5\epsilon_{1i}, 1 \leq i \leq 500$, where ϵ_{1i} is t distribution with degree of freedom 3; $D_i = 0.2 + 0.5\epsilon_{2i}, 501 \leq i \leq 1000$, where ϵ_{2i} is standard normal distribution.

(5) $D_i = 1 + \epsilon_i, 1 \leq i \leq 500$ and $D_i = 0.5 + \epsilon_i, 501 \leq i \leq 1000$, where ϵ_i is t distribution with degree of freedom 2.

In cases (1) and (4) the error term comes from two distinct families of distribution (normal and t distribution); in cases (1), (2), (3), and (4), the first 500 matched pairs have substantial treatment effects accompanied with large variability while the rest 500 matched pairs have small treatment effects and smaller variability in the error terms

Sampling Situation	Γ	Mean difference	M-statistic	Group M-statistic
Situation (1)	1	1.000	1.000	1.000
	2	0.997	0.853	0.996 (-0.001)
	3	0.145	0.002	0.133 (-0.012)
	4	0.000	0.000	0.000
Situation (2)	1	1.000	1.000	1.000
	2	1.000	1.000	1.000
	3	0.998	1.000	1.000
	4	0.769	0.848	0.926
	5	0.186	0.153	0.296
Situation (3)	1	1.000	1.000	1.000
	2	1.000	1.000	1.000
	3	1.000	0.998	1.000
	4	0.993	0.588	0.991 (-0.002)
	5	0.675	0.051	0.656 (-0.019)
Situation (4)	1	1.000	1.000	1.000
	2	1.000	1.000	1.000
	3	0.965	0.991	0.998
	4	0.504	0.346	0.708
	5	0.090	0.012	0.146
Situation (5)	1	1.000	1.000	1.000
	2	0.900	1.000	1.000
	3	0.238	0.843	0.798 (-0.045)
	4	0.016	0.113	0.093 (-0.020)
	5	0.000	0.003	0.003

Table 2: Simulated power of sensitivity analysis under different sampling situations for $I = 1000$ matched pairs. Each sampling situation is repeated 10,000 times. There are two binary covariates of interest and effect modification happens for the first 500 matched pairs. The details have been described in the main text. The power of group M-statistics in each sampling situation is in bold if it achieves the highest power (including tied results); otherwise, we show in parenthesis the difference between the power of group M-statistics and the statistics that achieves the highest power and the corresponding power of that test statistic is in bold.

and it is in these cases that trimming all the matched pairs together would potentially confuse effect modification with outliers and result in smaller sensitivity value; in case (5), we have moderate effect modification with the same t-distribution for all the 1000 matched pairs.

We note that the proposed group M-statistic performs well in all the settings as shown in Table 2. In particular, it achieves larger insensitivity to bias in sampling situation (2) and (4) where the error term follows a t-distribution in the subgroup with much larger treatment effect. If other statistics have higher power, then the differences are almost negligible whether it is achieved by using raw data as in situation (3) or using M-statistics in a conventional manner as in situation (5). Overall, the simulation results show that our group M-statistic that does subgroup-aware scoring works well across different settings.

6. Discussion

In this article, we introduced a novel scoring strategy that accounts for the subgroup structures when using M-statistics. By scoring the matched pairs within each subgroup, the proposed method works well with the malaria control data and recovers the lost insensitivity to bias if using the conventional M-statistics compared to using the raw data (which lacks robustness). We demonstrated that the proposed method is suitable in various settings and enjoys superior or non-inferior reported sensitivity value compared to using raw data and using conventional M-statistics.

Importantly, the practice of scoring all the matched pairs together when using conventional M-statistics would potentially confuse the outliers with effect modification. This practice does not distinguish whether the scores at the tail of the treated-minus-control differences come from the subgroup with a large treatment effect (e.g. people with age less than ten in the malaria control example) or from the randomness of long tail distribution of the error term. Figure B.1 shows the resulting distribution of treated-minus-control scores for different age groups given by three test statistics. We note that the contribution of pairs from people with age less than ten was trimmed away when using conventional M-statistics which resulted in inferior performance in sensitivity analysis. By contrast, the proposed group M-statistics preserved pairs with large treatment effect and Figure B.2 further shows that the outliers are properly trimmed in each interaction group by

our proposed method.

Moreover, it is worth mentioning that the meaning of the test and the sensitivity analysis conducted by using submax method is to view the population as a whole and test the global null hypothesis of no treatment effect (Lee et al., 2018). If we would like to make further inferences for different subgroups, then it is required that the scores used in a given component test should only depend on the matched pairs within these subgroups involved in the test, and our proposed method could be naturally used in such scenario while the conventional M-statistics would not be feasible (see further discussion in section 4 of Lee et al. (2018)). Meanwhile, as is the case in the malaria control data, the error term of different subgroups could be different since now we condition on a covariate. The interplay between the error distributions and the effect sizes across different subgroups together explains how sensitivity analysis behaves in the setting of effect modification (Lee et al., 2018; Hsu et al., 2013). We used the example of Huber’s M-statistics to illustrate the proposed subgroup-aware scoring approach and it is straightforward to extend this idea to other scoring procedures when we would like to make the scoring function dependent on the covariates. In general, if we want to leverage the effect modification or conduct the subgroup comparisons, it is critical to be aware of subgroup structure throughout the analysis procedure.

Acknowledgment

The authors thank Eric Sun for helpful discussions and preparation of an initial version of the figures.

Appendix A. R package

An R package `groupmscorev` will be available on CRAN.

Appendix B. Supplementary figures

Two supplementary figures are provided for visualizing the consequent distribution of treated-minus-control differences by using three test statistics on malaria control data.

References

- Gastwirth, J.L., Krieger, A.M., Rosenbaum, P.R., 2000. Asymptotic separability in sensitivity analysis. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)* 62, 545–555.
- Hansen, B.B., Rosenbaum, P.R., Small, D.S., 2014. Clustered treatment assignments and sensitivity to unmeasured biases in observational studies. *Journal of the American Statistical Association* 109, 133–144.
- Hsu, J.Y., Small, D.S., Rosenbaum, P.R., 2013. Effect modification and design sensitivity in observational studies. *Journal of the American Statistical Association* 108, 135–148.
- Lee, K., Small, D.S., Rosenbaum, P.R., 2018. A powerful approach to the study of moderate effect modification in observational studies. *Biometrics* 74, 1161–1170.
- Maritz, J.S., 1979. A note on exact robust confidence intervals for location. *Biometrika* 66, 163–166.
- Molineaux, L., Gramiccia, G., 1980. The Garki project: research on the epidemiology and control of malaria in the Sudan savanna of West Africa. World Health Organization, Geneva.
- Rosenbaum, P.R., 1987. Sensitivity analysis for certain permutation inferences in matched observational studies. *Biometrika* 74, 13–26.
- Rosenbaum, P.R., 2002. *Observational Studies*. Springer Series in Statistics. 2 ed., Springer New York, NY.
- Rosenbaum, P.R., 2007. Sensitivity analysis for m-estimates, tests, and confidence intervals in matched observational studies. *Biometrics* 63, 456–464.
- Rosenbaum, P.R., 2013. Impact of multiple matched controls on design sensitivity in observational studies. *Biometrics* 69, 118–127.
- Rosenbaum, P.R., 2017. submax: Effect Modification in Observational Studies Using the Submax Method. R package version 1.1.1.

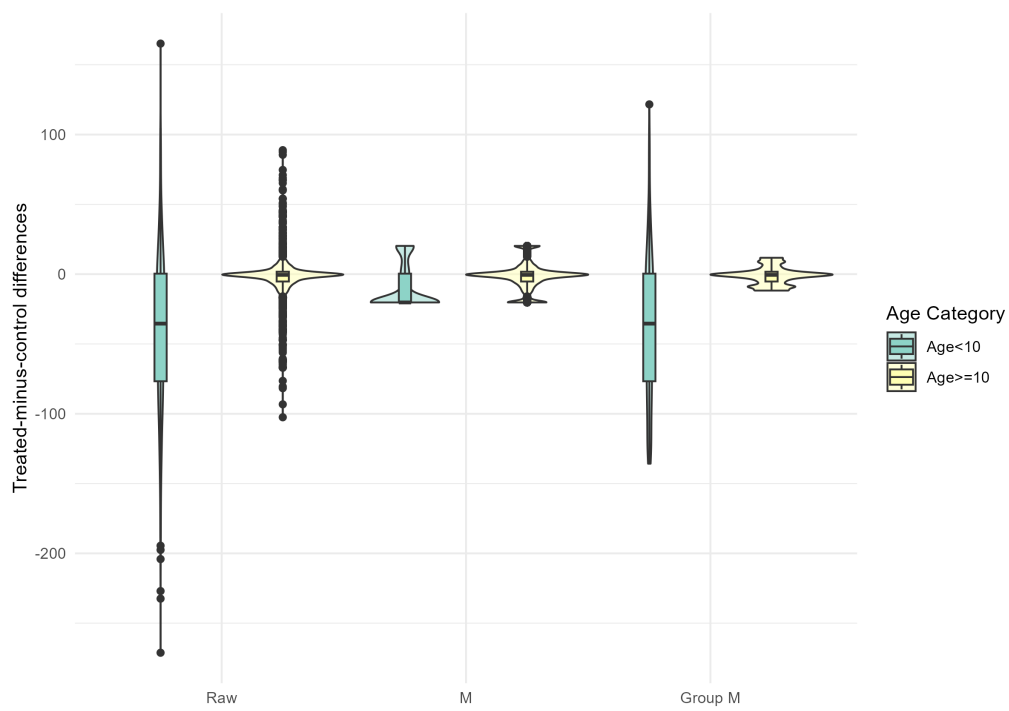


Figure B.1: Treated-minus-control differences across different statistics for people with age less than 10 and larger than 10 in the malaria control data. For the M-statistic, we multiply back the parameter h_0 in the figure to make the comparisons on the same scale with the raw data and the group M-statistic. It is essentially the same when we do inference.

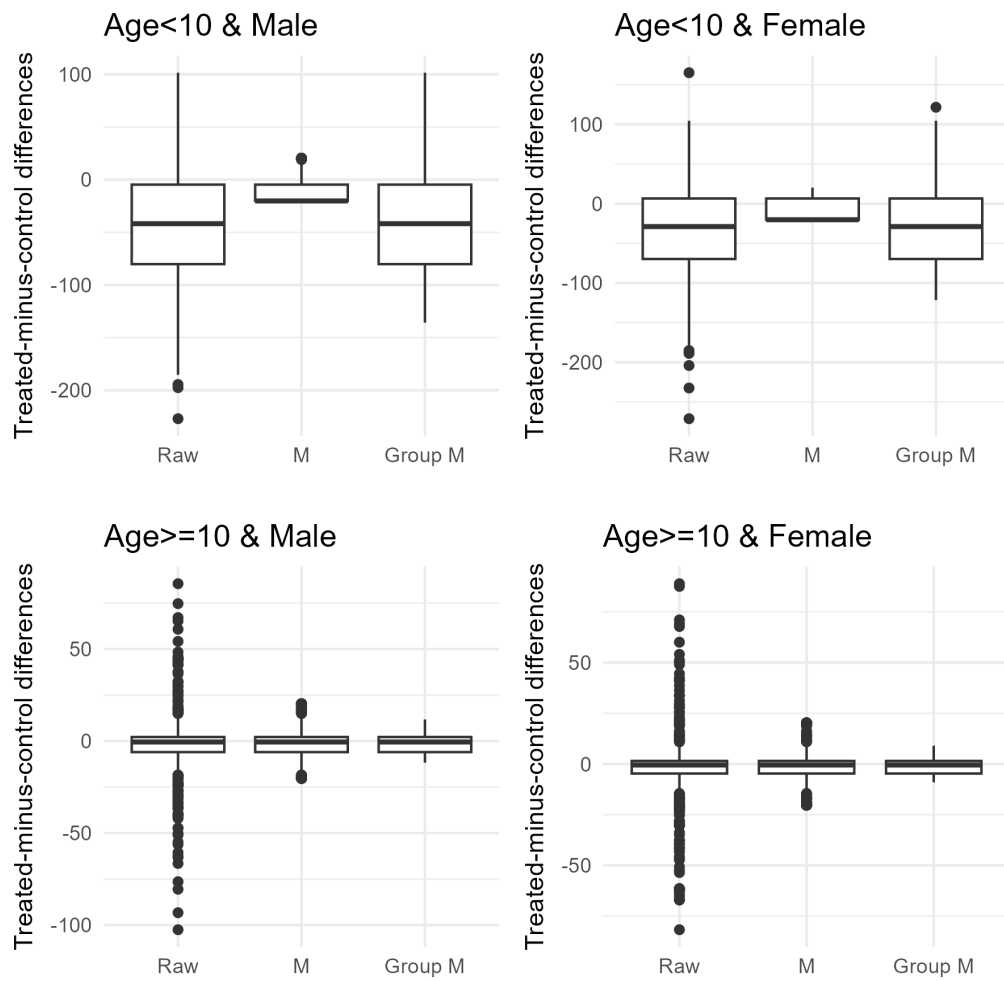


Figure B.2: Treated-minus-control differences across different statistics for people in four non-overlapped subgroups defined by age and gender. For the M-statistic, we multiply back the parameter h_0 in the figure to make the comparisons on the same scale with the raw data and the group M-statistic.